

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

Attention during natural vision warps semantic representation across the human brain.

### Permalink

<https://escholarship.org/uc/item/8k45n8rx>

### Journal

Nature neuroscience, 16(6)

### ISSN

1097-6256

### Authors

Çukur, Tolga  
Nishimoto, Shinji  
Huth, Alexander G  
et al.

### Publication Date

2013-06-01

### DOI

10.1038/nn.3381

Peer reviewed



Published in final edited form as:

*Nat Neurosci.* 2013 June ; 16(6): 763–770. doi:10.1038/nn.3381.

## Attention During Natural Vision Warps Semantic Representation Across the Human Brain

Tolga Çukur<sup>a</sup>, Shinji Nishimoto<sup>a,b</sup>, Alexander G. Huth<sup>a</sup>, and Jack L. Gallant<sup>a,c,d</sup>

<sup>a</sup>Helen Wills Neuroscience Institute, University of California, Berkeley, CA 94720, USA

<sup>c</sup>Program in Bioengineering, University of California, Berkeley, CA 94720, USA

<sup>d</sup>Department of Psychology, University of California, Berkeley, CA 94720, USA

### Abstract

Little is known about how attention changes the cortical representation of sensory information in humans. Based on neurophysiological evidence, we hypothesized that attention causes tuning changes to expand the representation of attended stimuli at the cost of unattended stimuli. To investigate this issue we used functional MRI (fMRI) to measure how semantic representation changes when searching for different object categories in natural movies. We find that many voxels across occipito-temporal and fronto-parietal cortex shift their tuning toward the attended category. These tuning shifts expand the representation of the attended category and of semantically-related but unattended categories, and compress the representation of categories semantically-dissimilar to the target. Attentional warping of semantic representation occurs even when the attended category is not present in the movie, thus the effect is not a target-detection artifact. These results suggest that attention dynamically alters visual representation to optimize processing of behaviorally relevant objects during natural vision.

### Keywords

Attention; Visual Search; Representation; Tuning Change; fMRI; Natural Stimuli; Movie

Attention is thought to increase information processing efficiency throughout the brain through several convergent mechanisms<sup>1</sup>. Neurophysiology studies in early visual areas have shown that spatial attention changes response baseline, response gain and contrast gain<sup>2–4</sup>. However, because the brain pools information across successive stages of processing, attentional modulation of baseline and gain at early stages likely causes changes

Users may view, print, copy, download and text and data- mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use: [http://www.nature.com/authors/editorial\\_policies/license.html#terms](http://www.nature.com/authors/editorial_policies/license.html#terms)

Correspondence should be addressed to: Jack L. Gallant, 3210 Tolman Hall #1650, University of California at Berkeley, Berkeley, CA 94720, <gallant@berkeley.edu>.

<sup>b</sup>Present address: Center for Information and Neural Networks, National Institute of Information and Communications Technology, Suita, Osaka, Japan

Supplementary Information is linked to the online version of this article.

**Author Contributions** T.Ç. and S.N. designed the experiments. T.Ç. and A.G.H. operated the scanner. T.Ç. conducted the experiments and analyzed the data. T.Ç. and J.L.G. wrote the manuscript. J.L.G. provided guidance on all aspects of the project.

The authors declare no competing financial interests.

in neuronal tuning in higher sensory and cognitive brain areas<sup>3,5</sup>. Indeed, feature-based attention can cause modest changes in tuning of single neurons even as early as V4<sup>5,6</sup>, but tuning changes of single neurons in pre-frontal cortex can be substantial<sup>7–9</sup>. Tuning shifts in single neurons change the way that information is represented across the neural population, warping the representation to favor certain signals at the expense of others<sup>5</sup>. Therefore, it has been proposed that tuning shifts reflect the operation of a matched-filter mechanism that optimizes task performance by expanding the cortical representation of attended targets<sup>5,10</sup>.

Attentional warping of cortical representation might be particularly valuable during demanding tasks such as natural visual search. Recent evidence from our laboratory suggests that the brain represents thousands of object categories by organizing them into a continuous semantic similarity space (Fig. 1a) that is mapped systematically across visual cortex<sup>11</sup>. Because natural scenes are cluttered with many different objects they may elicit patterns of brain activity that are widely distributed across this semantic space, making target detection difficult. Attention could dramatically increase sensitivity for the target and improve target detection under these demanding conditions<sup>5</sup>, by expanding the cortical representation of behaviorally relevant categories and compressing the representation of irrelevant categories (Fig. 1b–c).

It is currently unknown whether attention warps the cortical representation of sensory information in the human brain. To search for evidence for this complex attentional effect, we exploited the fact that attention would expand the representation of an attended category by causing neural populations throughout visual and non-visual cortex to shift tuning toward the target<sup>5–9</sup> (see Supplementary Fig. 1 for alternative hypotheses). We thus hypothesized that visual search for a single object category should cause tuning shifts in single voxels measured by functional MRI (Fig. 1d–f).

To identify semantic tuning shifts we measured category tuning in single voxels during a natural category-based visual search task (Fig. 2). We recorded whole-brain fMRI data from five human subjects while they viewed 60 minutes of natural movies (see Online Methods). Subjects maintained steady fixation while covertly searching for ‘humans’ or ‘vehicles’. These categories were used because they are quite distinct from one another, they occur commonly in real-world scenes, and they are common targets of visual search<sup>12,13</sup>.

Category-based attention tasks have been used in several previous fMRI experiments<sup>12,14,15</sup>. However, these earlier studies used a small set of object categories and region-based data analysis procedures. Therefore, they did not explore voxel-based tuning and could not distinguish voxel-based changes in tuning from changes in response baseline or gain. To maximize our ability to detect tuning changes in single voxels, we used complex natural movie stimuli containing hundreds of different object and action categories<sup>16,17</sup>. To remove attentional effects on response baseline and gain, we normalized the blood-oxygen-level-dependent (BOLD) responses of each voxel to have zero mean and unit variance individually within each attention condition before further modeling. This procedure allowed us to clearly separate tuning changes from simple modulation of response baseline or gain.

We then employed a voxel-wise modeling approach developed previously in our laboratory to obtain accurate estimates of category tuning in single cortical voxels, and in each individual subject<sup>11,18–21</sup>. The WordNet lexicon<sup>22</sup> was used to label 935 object and action categories in the movies (Supplementary Fig. 2). Regularized linear regression was used to fit voxel-wise models that optimally predicted the measured BOLD responses from the categorical indicator variables (Supplementary Fig. 3). Separate models were estimated using data acquired during visual search for ‘humans’ and for ‘vehicles’. The resulting model weights give the category tuning vectors for each voxel, under each attention condition.

## Results

Attentional changes in semantic representation can be inferred by comparing category tuning vectors across attention conditions (see Fig. 3 for tuning vectors for one voxel located in lateral occipital complex). However, inferences drawn from this comparison will only be justified and functionally important if the fit category models can successfully predict BOLD responses to novel natural stimuli. To address this issue we validated the prediction performance of category models on separate data reserved for this purpose. Prediction scores were defined as the Pearson’s correlation between the BOLD responses measured in the validation dataset, and those predicted by the fit models (see Online Methods). All statistical significance levels were corrected for multiple comparisons using false-discovery-rate control<sup>23</sup>.

We find that category models provide accurate response predictions across many regions of visual and non-visual cortex (Supplementary Fig. 4). Overall,  $83.7 \pm 5.12\%$  (mean  $\pm$  s.d. across subjects) of cortical voxels are significantly predicted by the category model (t-test,  $p < 0.05$ ). The category model explains more than 20% of the response variance in  $11.60 \pm 5.84\%$  (mean  $\pm$  s.d.) of these voxels across subjects. These results suggest that category tuning vectors accurately reflect category responses of many cortical voxels during visual search.

If attentional tuning changes are statistically significant, then category models for individual attention conditions should yield better response predictions than a null model fit by pooling data across conditions. To assess significance we therefore compared the prediction scores obtained from category models to those obtained using null models. We find that  $59.57 \pm 8.31\%$  (mean  $\pm$  s.d. across subjects) of cortical voxels exhibit significant tuning changes (t-test,  $p < 0.05$ ). Across subjects  $17.13 \pm 0.97\%$  (mean  $\pm$  s.d.) of these voxels also have high prediction scores (above mean plus 1 standard deviation), yielding 4245–7785 well-modeled voxels in individual subjects. Because all responses were z-scored individually within attention conditions, these results cannot be explained by additive or multiplicative modulations of responses in single voxels. Therefore, they indicate that attention causes significant tuning changes in many cortical voxels.

Our experiment used a category-based attention task that required attention to ‘humans’ or ‘vehicles’. However, complex natural movies may contain low-level features that are correlated with these semantic categories. Do the attentional tuning shifts shown here reflect

category-based attention or rather are they due to attention to correlated low-level features? To address this issue we fit simpler structural encoding models that reflect tuning for elementary features such as spatio-temporal frequency, orientation, and eccentricity (see Online Methods for details). We then compared predictions of category models and structural models across well-modeled voxels that showed significant tuning shifts in the category-based attention task.

We find that the average prediction score of structural models is only  $0.22 \pm 0.03$  (mean  $\pm$  s.d. across subjects), which is significantly lower than that of category models ( $0.54 \pm 0.11$ , randomized t-test,  $p < 10^{-4}$ ). We also find that the percentage of response variance explained by structural tuning shifts is only  $1.53 \pm 0.68\%$  (mean  $\pm$  s.d.), which is significantly lower than that explained by category-based tuning shifts ( $13.57 \pm 7.65\%$ , Wilcoxon signed-rank test,  $p < 10^{-4}$ ). These findings suggest that tuning for elementary visual features cannot account for category-based tuning shifts measured here.

Next, we asked whether these changes in category tuning are consistent with the tuning-shift hypothesis<sup>5</sup>, in which attention warps semantic representation to favor behaviorally relevant categories at the expense of irrelevant categories. The tuning-shift hypothesis makes three explicit and diagnostic predictions about how attention alters semantic representation. First, it predicts that attention causes tuning shifts toward the attended category when the targets are present, expanding the representation of the attended category. Second, it predicts that attention causes tuning shifts toward the attended category even when no targets are present. Finally, it predicts that attention expands the representation of unattended categories that are semantically similar to the target, and compresses the representation of categories that are semantically dissimilar to the target. We tested the tuning-shift hypothesis by evaluating each of these predictions in turn.

### Tuning Shifts in the Presence of Targets

To determine whether attention causes tuning shifts toward the attended category when the targets are present, we first projected voxel-wise tuning vectors measured during visual search into a continuous semantic space. The semantic space was derived from principal components analysis of tuning vectors measured during a separate passive-viewing task (see Online Methods). Different voxels that are tuned for semantically similar categories will project to nearby points in this space. We then visualized the distribution of tuning across well-modeled voxels that have significant category models (t-test,  $p < 0.05$ ). We find that most well-modeled voxels are selectively tuned for the attended category, and attention causes tuning shifts in most of these voxels (Fig. 4a).

We quantified the magnitude and direction of tuning shifts across attention conditions by measuring the selectivity of voxel tuning for ‘humans’ or ‘vehicles’ under each condition (see Online Methods; Supplementary Figs. 5 and 6a–e). We then computed a tuning shift index (TSI) that summarizes the difference in selectivity for the attended versus unattended category (see Online Methods). Under this scheme a voxel that shifts toward the attended category will have a positive TSI. We find that the mean TSI across well-modeled voxels is significantly greater than 0 in all subjects (Wilcoxon signed-rank test,  $p < 10^{-6}$ ; Supplementary Fig. 7). Because all responses were z-scored individually within each

attention condition before TSI values were calculated, these tuning shifts cannot be explained by changes in voxel response baseline or gain (see also Discussion). Thus, these results are consistent with the view that attention changes tuning to expand the representation of the attended category.

### Cortical Distribution of Tuning Shifts

Previous neurophysiology studies suggest that tuning shifts should be widespread across the brain, extending from higher-order visual areas into frontal cortex<sup>5–9</sup>. To visualize the distribution of tuning shifts across cortex, we projected TSI values onto cortical flat maps. We find that voxels in many different brain regions shift their tuning toward the attended category (Fig. 4b, see also Supplementary Fig. 6a–e). (Interested readers may explore the datasets at <http://gallantlab.org/brainviewer/cukuretal2013>.) These include most of ventral-temporal cortex; the lateral-occipital and intraparietal sulci; the inferior and superior frontal sulci; and the dorsal bank of the cingulate sulcus (see also Supplementary Figs. 9 and 10). In contrast to most brain regions, voxels in the precuneus, temporal-parietal junction, anterior prefrontal cortex, and anterior cingulate sulcus shift their tuning away from the attended category. This finding suggests that these brain areas are involved in distractor detection and in error monitoring during visual search<sup>24,25</sup>.

To examine how specific brain areas change their representations of attended and unattended categories, we performed detailed analyses of tuning shifts in several common regions-of-interest. We find that regions in higher-order visual cortex and more anterior brain areas have high prediction scores, indicating that tuning shifts in these regions are functionally important (Fig. 5a). TSI is small in retinotopic early visual areas, but it is significantly larger in more anterior brain areas that correspond to later stages of visual processing (Wilcoxon signed-rank test,  $p < 10^{-6}$ ; Fig. 5b). This result implies that attentional tuning shifts become progressively stronger towards later stages of processing. We also find that these tuning shifts occur for both attended (i.e., ‘humans’ and ‘vehicles’; Fig. 5c) and unattended categories (Wilcoxon signed-rank test,  $p < 10^{-6}$ ; Fig. 5d). This finding is consistent with an attentional mechanism that alters the representation of the entire semantic space during visual search (see Supplementary Fig. 1d). Finally, we find that tuning changes for attended categories account for a relatively larger fraction of the overall tuning change in more anterior brain areas compared to earlier visual areas (Fig. 5c), while those for unattended categories account for a relatively smaller fraction of tuning changes (Fig. 5d). Taken together, these results suggest that more anterior brain areas are primarily involved in representing the attended category, and that visual representations in more frontal areas are relatively more dependent on the search task than are those at earlier stages of visual processing<sup>5–9,26</sup>.

### Tuning Shifts in the Absence of Targets

The second prediction of the tuning-shift hypothesis is that attention causes tuning changes even when no targets are present. To address this issue we estimated voxel tuning using only those segments of the movies that did not contain ‘humans’ or ‘vehicles’. (Note also that because any systematic differences in arousal, respiration, and spatial attention across attention conditions are most likely to occur when the targets are present, this analysis also

serves as a powerful control against such nuisance factors; see Online Methods for additional controls). Because data recorded when the targets were present were excluded from analysis, tuning for the attended categories cannot be assessed directly. However, our modeling framework allows us to measure tuning shifts for the remaining categories, and to infer the direction of shifts with respect to the attended categories from these measurements.

To assess the direction of tuning shifts in the absence of the targets, we projected the tuning vectors estimated in the absence of the targets into the semantic space. We find that voxels in many brain regions shift their tuning toward the attended category even when no targets are present (Fig. 6, Supplementary Fig. 11a–e; explore the datasets at <http://gallantlab.org/brainviewer/cukuretal2013>). The mean TSI across the population of well-modeled voxels is significantly greater than 0 in all subjects (Wilcoxon signed-rank test,  $p < 10^{-6}$ ; Supplementary Fig. 8). These results demonstrate that attention causes tuning shifts toward the attended category even when no targets are present, and that attentional tuning shifts are not a mere consequence of target detection.

### Semantic Representation of Unattended Categories

The third prediction of the tuning-shift hypothesis is that attention expands the representation of categories that are semantically similar to the attended category, even when no targets are present. If the representation of an unattended category is expanded, its representation should shift toward the representation of the attended category (i.e., the region of the semantic space that many voxels are tuned for). To address this issue we assessed how the similarity between representations of unattended and attended categories changed across attention conditions. The similarity between representations of two categories was measured using Pearson's correlation between corresponding BOLD-response patterns across well-modeled voxels<sup>27</sup>. Responses for unattended and attended categories were estimated using target-absent and target-present movie segments, respectively. We find that during search for 'humans', representations of animals, body parts, action verbs, and natural materials shift toward the representation of 'humans'. During search for 'vehicles', representations of tools, devices, and structures shift toward the representation of 'vehicles' (Wilcoxon signed-rank test,  $p < 10^{-4}$ ; Fig. 7). This result suggests that attention expands the representation of unattended categories that are semantically similar to the target, at the expense of semantically dissimilar categories.

### Discussion

Our results indicate that category-based attention during natural vision causes semantic tuning changes that cannot be explained by additive or multiplicative response modulations in single voxels. These tuning changes alter the cortical representation of both attended and unattended categories. Furthermore, attentional changes in tuning for unattended categories occur even when the attended categories are not present in the movie. These effects are consistent with an attentional mechanism that acts to expand the representation of semantic categories nearby the target in the semantic space at the cost of compressing the representation of distant categories.



Because this study measured hemodynamic changes we cannot make direct inferences about the underlying neural mechanisms mediating tuning shifts. Several possible neural mechanisms might conceivably contribute to semantic tuning changes in single voxels. When the targets are present in the display, then it is possible that changes in response baseline or gain of single neurons that are tuned to the attended targets contribute to tuning changes. However, tuning changes for unattended categories observed when no targets are present cannot be explained by this mechanism: because the attended categories were never present in these cases, neurons tuned only to the attended categories never entered into the model estimation procedure and therefore they could not have any effect on estimated voxel-wise tuning curves.

Our results are consistent with existing neurophysiology studies that have demonstrated tuning shifts in single neurons in as early as area V4<sup>5,6</sup>, and which have shown far stronger tuning shifts at relatively higher levels of visual and cognitive processing<sup>7–9</sup>. Some of these single neuron studies have reported that tuning shifts are consistent with a matched-filter mechanism that shifts tuning toward the attended target, expanding the representation of attended stimuli at the cost of unattended stimuli. Our results are also consistent with theoretical expectations based on the anatomical structure of the cortical hierarchy: because neurons pool information across successive stages of processing, attentional modulation of baseline or gain at one level must inevitably cause tuning changes at subsequent levels<sup>3,5</sup>. Thus, it is reasonable to expect that changes in voxel tuning at least partly reflect tuning shifts in individual neurons within the underlying neural population.

Although natural movies have strong face validity, correlations inherent in natural movies could potentially complicate interpretation of the results. We took several measures to ensure that stimulus correlations did not confound our results. First, the collection of movies used in the experiments was highly diverse. Second, we used a regression-based modeling approach that minimizes the effect of residual correlations on the fit models. Third, we performed control analyses on raw BOLD responses to rule out biases due to correlations between attended and unattended categories (see Online Methods).

Given that our data are finite, there is always some chance that residual correlations may introduce some bias in the results. However, artificial stimuli that contain only a small number of categories introduce much more substantial and pernicious bias, and so are more likely to lead to misinterpretation. Interpretation of experiments that use limited stimulus sets inevitably rely on a strong assumption of linearity, that is, that responses to multiple objects in a natural context will be predictable from responses to isolated objects. In contrast, natural stimuli do not require any such linearity assumptions. Note, however, that this important issue is really not relevant to this study. The main goal of this study is not to measure tuning, but rather to measure changes in tuning between different search tasks. Because natural stimuli have high ecological relevance for natural visual search, natural movies appear to be better suited for these measurements.

An important question to be answered is the role of bottom-up processing versus top-down feedback in measured tuning changes. Because we used the same movie stimulus for the two separate search tasks in our experiment, all attentional tuning changes between the two tasks



must necessarily reflect top-down modulatory effects. We find small tuning shifts in retinotopic early visual areas, but significantly larger tuning shifts in higher visual areas in occipito-temporal cortex and relatively more anterior brain areas. We also show that tuning shifts cannot be explained by response modulations for lower-level visual features that are known to be represented in early visual areas. These results imply that attentional modulations primarily warp semantic representation at later stages of visual processing. However, the slow nature of BOLD responses makes it difficult for any fMRI study to measure the temporal relationship between signals arising in different brain areas at these later stages of processing. We plan to investigate this issue in the future with neurophysiology studies in animals to achieve sufficiently high temporal resolution.

The way that attention optimizes target detection depends not only on the target, but also on the similarity between the target and the distractors<sup>28</sup>. If the target is very different from the distractors, then target detection can be optimized by shifting tuning toward the target<sup>5</sup>. However, if the target is very similar to the distractors, target detection can be improved by enhancing the representation of task-irrelevant features that optimally distinguish the target from the distractors<sup>29</sup>. In the study reported here the attentional targets were highly distinct ('humans' and 'vehicles'), so it is natural to expect that tuning should shift toward the target. An important topic for future research will be to determine whether attention causes tuning shifts toward task-irrelevant features when the target and distractors are very similar.

In conclusion, we find that natural visual search for a single category warps the entire semantic space, expanding the representation of nearby semantic categories at the cost of more distant categories. This effect suggests a more dynamic view of attention than is assumed under the conventional view that attention is a simple mechanism that merely modulates the baseline or gain of labeled lines. This dynamic mechanism can improve the effective resolution of the visual system for natural visual search, and it likely enables the use of limited neural resources to perform efficient search for many different object categories. Overall, these findings help explain the astounding human ability to perform complex visual tasks in an ever-changing natural environment.

## Online Methods

### Subjects

Five healthy adult volunteers (five males) with normal or corrected to normal vision participated in this study: S1 (age 30), S2 (age 32), S3 (age 25), S4 (age 25), and S5 (age 26). The experimental procedures were approved by the Institutional Review Board at the University of California, Berkeley (UCB); and written informed consent was obtained from all subjects.

### Stimuli

For each attention condition in the main experiment, 1800 sec of continuous color natural movies (24°x24°, 512x512 pixels) were presented without repetition in a single session. The stimuli were compiled by combining many short clips (10–20 sec) from a diverse selection of natural movies<sup>19</sup>. Only 'humans' or only 'vehicles' were present in 450 sec both, the two categories co-occurred in 450 sec, and both categories were absent during 450 sec.

'Humans' and 'vehicles' appeared in highly diverse scenes and in many different positions, sizes and viewpoints. A fixation spot ( $0.16^\circ$  square) was superimposed on the movies and its color was alternated at 1 Hz, rendering it continuously visible. The stimuli were presented at a rate of 15 Hz using an MR-safe projector (Avotec Inc., Stuart, FL) and a custom-built mirror system.

### Experimental Paradigm

Each subject participated in a total of seven scan sessions. Functional localizer, retinotopic mapping and anatomical data were collected in two sessions. Functional scans for the main experiment were collected in a single scan session. To increase sensitivity for the analysis performed in the absence of the target stimuli, another session of functional data was collected using the same experimental design, but with a different set of movie clips. To construct the continuous semantic space, 7200 sec of natural movies were presented in three separate sessions while subjects performed a passive-viewing task.

In the main experiment, subjects fixated continuously while covertly searching for 'humans' or 'vehicles' in natural movies. To ensure continuous vigilance subjects depressed the response button continuously whenever an exemplar of the attended category was present in the movies. The data for each attention condition were recorded within 3 separate 10-min runs. The movie clips within each run were selected randomly without repetition. To avoid sampling bias, an identical set of movie clips were presented for both attention conditions. The presentation order of these clips was counterbalanced across the conditions. Four mutually exclusive classes of stimuli (i.e., only 'humans', only 'vehicles', both 'humans' and 'vehicles', and neither 'humans' nor 'vehicles') were randomly interleaved and evenly distributed within and across the runs. The attended category was fixed within each run. The attention conditions were alternated in consecutive runs. A cue word, 'humans' or 'vehicles', was displayed prior to each run to indicate the attended category. To compensate for hemodynamic transients caused by movie onset, each run was preceded by the last 10 sec of that run. Data collected during the transient period were discarded.

### MRI Protocols

MRI data were acquired on a 3 T Siemens scanner located at the University of California, Berkeley using a 32-channel head coil. Functional data were acquired using a  $T_2^*$ -weighted gradient-echo EPI sequence customized with a water-excitation radiofrequency pulse to prevent contamination from fat signal. The following parameters were prescribed: repetition time = 2 sec, echo time = 34 msec, flip angle =  $74^\circ$ , voxel size =  $2.24 \times 2.24 \times 3.5$  mm<sup>3</sup>, field-of-view =  $224 \times 224$  mm<sup>2</sup>, and 32 axial slices to cover the entire cortex. Head motion was minimized with foam padding. To reconstruct cortical surfaces, anatomical data were collected with  $1 \times 1 \times 1$  mm<sup>3</sup> voxel size and  $256 \times 212 \times 256$  mm<sup>3</sup> field-of-view using a three-dimensional  $T_1$ -weighted MP-RAGE sequence. The anatomical and retinotopic mapping data for subjects S2 and S3 were obtained on a 1.5 T Philips Eclipse (Philips Medical Systems, NA, Bothell, WA) scanner.

## Data Pre-processing

Functional scans were intra- and inter-run aligned using the Statistical Parameter Mapping toolbox (SPM8, <http://www.fil.ion.ucl.ac.uk/spm/software/spm8/>). All volumes were aligned to the first image from the first functional run for each subject. Non-brain tissue was excluded from further analysis using the Brain Extraction Tool (BET, <http://www.fmrib.ox.ac.uk/analysis/research/bet/>). Voxels whose BOLD responses are primarily driven by button presses were identified using a motor localizer that included a button-press task. The identified voxels were contained within the primary motor, somatosensory motor and premotor cortices, and voxels within these regions were excluded from analysis. The cortical surface of each subject was reconstructed from anatomical data using Caret5 (<http://www.nitrc.org/projects/caret/>). Cortical voxels were identified as the set of voxels within a 4-mm radius of the cortical surface. Subsequent analyses were restricted to 47125-53957 cortical voxels identified for the various subjects.

The low-frequency drifts in voxel responses were estimated using a 240-sec-long cubic Savitzky-Golay filter for each run (10 min). The drifts were removed from the responses, which were then normalized to have zero mean and unit variance. Neither spatial nor temporal averaging was performed on the data during pre-processing and model-fitting stages. The data from separate subjects were not transformed into a standard brain space.

Functional localizer and retinotopic mapping data were used to assign voxels to the corresponding regions-of-interest (ROI)<sup>11</sup>. All functional ROIs were defined based on relative response levels to contrasting stimuli (t-test,  $p < 10^{-5}$ , uncorrected).

## Category Model

The object and action categories in each one-second clip of the natural movie stimulus were manually labeled using terms from the WordNet lexicon<sup>22</sup>. Three naïve raters performed the labeling, and potential conflicts were resolved by conferral among all raters. In WordNet, words are grouped into sets of synonyms according to the concepts they describe, and are organized into a hierarchical network of semantic relations based on word meaning. By definition, the existence of a category in a given scene indicates the existence of all of its superordinate categories. For example, if a clip is labeled with ‘child’, it also contains the following categories: ‘offspring’, ‘relative’, ‘person’, ‘organism’, ‘living thing’, ‘whole’, ‘object’, and ‘entity’. To facilitate labeling, the raters exploited these hierarchical relationships in WordNet. The raters initially labeled 604 object and action categories, and inferred the presence of 331 superordinate categories from these initial labels.

A stimulus time-course (categories x seconds) was then formed by using a binary variable to indicate the presence or absence of each category in each one-second movie clip. The category model fit to each voxel describes evoked responses as a weighted linear combination of these indicator variables. The predicted response of each voxel to any category is the sum of weights for all the categories it encompasses (including itself). In other words, the weight for each category is the estimated difference between the response to that specific category and the cumulative response to all of its superordinate categories.

Retinotopically-organized early visual areas (V1–V4) are selective to structural characteristics of visual stimuli<sup>19</sup>. To ensure that model fits were not biased by structural differences in movie clips, one additional regressor was included in the model that characterizes the total motion energy in each one-second clip. This regressor was computed as the average response of 2139 space-time quadrature Gabor filter pairs to the movie stimuli. The filters were selected to cover the entire image space (24°x24°), and reflected a wide range of preferred receptive-field sizes, orientations, and spatiotemporal frequencies. In addition, to ensure that semantic tuning changes do not simply reflect tuning changes for elementary visual features, a separate structural model (with all 2139 filter pairs) was fit to each voxel.

To ensure that the results were not biased by the hierarchical relationships in WordNet, reduced category models were fit using the subset of regressors for the 604 initially-labeled categories. The data presented in this study was also analyzed within this separate framework, and no significant discrepancies were observed in the obtained results. Furthermore, the original full category model outperformed this reduced category model in terms of prediction accuracy of BOLD responses (see Supplementary Fig. 12). This indicates that the full category model provides a better description of category selectivity in cortical voxels.

### Model Fitting

The model for each attention condition was fit separately to 1800 sec of stimuli and responses. The stimulus time-course was down-sampled by a factor of 2 to match the sampling rate of the measured BOLD responses. To model the slow hemodynamic response, each category was assigned a distinct time-inseparable finite impulse response filter with delays restricted to 2–6 sec prior to the BOLD responses. All model parameters were simultaneously fit using L2-regularized linear regression.

To assess the significance of attentional tuning changes, a jackknifed model training/validation procedure was repeated 1000 times. At each turn, 20% of the samples were randomly held out to validate the model performance. The regularization parameter ( $\lambda$ ) for regression was selected with 10-fold cross-validation on the remaining 80% of training samples. These samples were further split into 10% testing and 90% training sets at each fold. The trained models were tested on the 10% held-out sets by computing prediction scores. Prediction score was taken as the correlation coefficient (Pearson's  $r$ ) between the actual and predicted BOLD responses. The optimal  $\lambda$  was determined for each voxel by maximizing the average prediction score. To prevent potential bias in the models, a final  $\lambda$ -value was selected as an intermediate between the optima for models from two attention conditions. The model-fitting procedures were performed with in-house software written in Matlab (The Mathworks, Natick, MA).

### Characterizing Tuning Shifts

Attentional tuning shifts toward the target will increase the degree of tuning selectivity – tuning strength – for the attended category. Therefore, the magnitude and direction of tuning shifts can be assessed by measuring the tuning strengths for ‘humans’ and ‘vehicles’

separately during each attention condition. Tuning strengths for ‘humans’ and ‘vehicles’ were quantified as the similarity between voxel tuning and idealized templates tuned solely for ‘humans’ and ‘vehicles’, respectively. The templates were constructed by identifying the set of labels that belong to these categories (Supplementary Fig. 5). Tuning strength for each category was then quantified as Pearson’s correlation between voxel tuning and the corresponding template.

$$s_{i,H} = \text{corr}(w_i, t_H)$$

$$s_{i,V} = \text{corr}(w_i, t_V)$$

Here,  $s_{i,H}$  is the tuning strength for ‘humans’, and  $s_{i,V}$  is the tuning strength for ‘vehicles’ during attention condition  $i$  ( $i=H$ : search for humans, and  $i=V$ : search for vehicles). Meanwhile,  $w_i$  is the voxel-wise tuning vector during condition  $i$ ; and  $t_H$  and  $t_V$  are the templates for ‘humans’ and ‘vehicles’, respectively.

Finally, a tuning shift index (TSI) was quantified using the measured tuning strengths for each voxel.

$$TSI = \frac{(s_{H,H} - s_{H,V}) + (s_{V,V} - s_{V,H})}{2 - \text{sign}(s_{H,H} - s_{H,V}) \cdot s_{H,V} - \text{sign}(s_{V,V} - s_{V,H}) \cdot s_{V,H}}$$

Here, the numerator measures the difference in tuning strength for the attended versus unattended category, summed across two attention conditions. Meanwhile, the denominator scales the TSI to range in  $[-1, 1]$ . Tuning shifts toward the attended category will yield positive TSIs, with a value of 1 in the case of a perfect match between voxel tuning and idealized template for the attended category. In contrast, tuning shifts away from the attended category will yield negative TSIs, with a value of  $-1$  in the case of a total mismatch between voxel tuning and idealized template for the attended category. Finally, a TSI of zero indicates that the voxel tuning did not shift between the two attention conditions.

Complementary tuning-shift analyses were performed in individual ROIs. For each attention condition, the mean tuning shift in each ROI was computed by averaging the TSI values of the corresponding set of voxels with significant models (t-test,  $p < 0.05$ , FDR corrected) and positive prediction scores.

## Eye-movement and Behavioral Controls

Eye movements are a legitimate concern in many experiments on visual perception and attention, especially when naïve subjects are tested. However, four lines of evidence demonstrate that eye movements were not a problem in our experiment, and that they could not have accounted for our results. First, all of the subjects tested in this experiment were highly trained psychophysical observers who had extensive experience in fixation tasks. Based on our previous work with trained and naïve subjects, we fully expect that our trained observers fixate much better than do the naïve subjects used in many attention experiments.

Second, there is no statistical evidence that fixation differs across attention conditions in the main and control analyses, for any of the observers. Subjects' eye positions were monitored at 60 Hz throughout the scans using a custom-built camera system equipped with an infrared source (Avotec Inc., Stuart, FL) and the ViewPoint EyeTracker software suite (Arrington Research, Scottsdale, AZ). The eye tracker was calibrated prior to each run of data acquisition. A nonparametric ANOVA test was used to determine systematic differences in the distribution of eye positions. The eye position distributions are not affected by attention condition ( $p>0.24$ ), or by target presence/absence ( $p>0.61$ ). To determine whether the results were biased by explicit eye movements during target or distractor detection, we also analyzed the distribution of eye positions during 250-msec, 500-msec and 1-sec windows around target onset, and target offset. The eye position distributions are not affected by target onset ( $p>0.26$ ) or offset ( $p>0.49$ ). Furthermore, there are no significant interactions between any of the aforementioned factors ( $p>0.14$ ). To determine whether the results were biased by rapid moment-to-moment variations in eye position, we examined the moving-average standard deviation of eye position within a 200-msec window (to capture potential saccades). There are no effects of attention condition ( $p>0.13$ ), target presence/absence ( $p>0.52$ ), target onset ( $p>0.47$ ) or target offset ( $p>0.17$ ), and there are no significant interactions between these factors ( $p>0.22$ ).

Third, while there may be some micro-saccade scale eye movements during covert visual search, there is no statistical evidence for a bias in the recorded BOLD responses across attention conditions. Specifically, there are no significant differences in BOLD responses due to interactions between the search task and scenes likely to contain the attended category or scenes that contain objects that share visual features with the attended category (two-way ANOVA,  $F<1.8$ ,  $p>0.18$ , FDR corrected). Because we measure attentional tuning shifts using BOLD responses, this analysis indicates that small eye movements could not have accounted for our results.

Finally, to further ensure that the results were not confounded by eye movements, we regressed the moving-average standard deviation of eye position out of the BOLD responses, and then we repeated the entire modeling procedure on these filtered data. Including this nuisance regressor did not affect the model fits or the results in any brain regions where the category model provided significant response predictions.

Behavioral responses were also recorded during the scans with a fiber-optic response pad (Current Designs Inc., Philadelphia, PA). A hit was defined as a button response detected within 1-sec of the target onset in the movies. A false alarm was defined as a button response when the target was absent from the movies. The behavioral performance, as measured by the sensitivity index ( $d'$ ), was compared across the two attention conditions using Wilcoxon rank-sum tests. Participants performed equally well when searching for either category, indicating that the task difficulty was balanced across attention conditions (Supplementary Fig. 13).

### Head-motion and Physiological-noise Controls

To ensure that our results were not biased by head motion or physiological noise, we used estimates of these nuisance factors to regress them out of the BOLD responses, and then we

repeated the entire modeling procedure on these filtered data. The moment-to-moment variations in head position were estimated during motion correction pre-processing. These six-parameter affine transformation estimates of head position were used to create head-motion regressors. The cardiac and respiratory states were recorded using a pulse oximeter and a pneumatic belt. These recordings were used to create pulse-oximetry and respiratory regressors as low-order Fourier series expansions of the cardiac and respiratory phases. The inclusion of these various nuisance regressors did not affect the model fits or the results in any brain region where the category model provided significant response predictions.

### Spatial-attention Controls

Given the stimulus correlations inherent in natural movies, differences in spatial attention across attention conditions might have confounded our results, even in the absence of targets. We performed two additional control analyses to ensure that the results derived from target-absent movie clips were not biased by stimulus correlations. First, all target-absent movie clips were coded to indicate whether they contained objects that shared visual features with ‘humans’ (i.e., scenes that contain animals, body parts, or animate motion) or with ‘vehicles’ (i.e., scenes that contain inanimate objects such as artifacts, buildings, or devices). Thereafter, a two-way ANOVA was performed on the evoked BOLD responses to determine whether there is any interaction between scene content and attended category. There are no significant interactions between scene content and attended category ( $F < 1.8$ ,  $p > 0.18$ , FDR corrected).

Second, all target-absent movie clips were coded to indicate whether humans were likely to appear (i.e., scenes that contain animate motion, tools for human use, buildings, or rooms) or whether vehicles were likely to appear (i.e., scenes of urban areas or cities, and scenes containing roads or highways). Another two-way ANOVA was performed on the evoked BOLD responses to determine whether there is any interaction between scene type and attended category. There are no significant interactions between scene type and attended category ( $F < 2.0$ ,  $p > 0.16$ , FDR corrected). Thus, there is no evidence for an interaction between scene content or type and the attentional target. These results suggest that the tuning shifts reported here are not biased by systematic differences in spatial attention across attention conditions.

### Construction of the Semantic Space

To construct the continuous semantic space, functional data were collected while subjects passively viewed 7200 sec of natural movies. Voxel-wise tuning vectors were estimated using these data and following identical procedures to the main experiment. A semantic space of cortical representation was then derived using principal components analysis (PCA) across the tuning vectors of cortical voxels (following procedures described in ref. 11). PCA ensures that voxels tuned for similarly represented categories project to nearby points in the semantic space, whereas voxels tuned for dissimilarly represented categories project to distant points. Each PC represents a distinct dimension of the semantic space, ordered according to percentage of variance explained. To maximize the quality of the semantic space, only the first six PCs were selected that captured approximately 30% of the variance. To perform analyses of attentional tuning changes in the semantic space, voxel-wise tuning



vectors obtained under different attention conditions were first projected onto these PCs. The results did not significantly vary with the number of PCs used to define the semantic space.

### Control Analysis in the Absence of Target Stimuli

A control analysis was performed to assess tuning changes for unattended categories in the absence of target stimuli. To increase sensitivity, additional functional data were collected in all subjects using the same experimental paradigm. A total of 1800 sec of stimuli were compiled from a different selection of movie clips than those used in the main experiment. To estimate tuning, the BOLD responses to movie clips where no ‘humans’ or ‘vehicles’ appeared were pooled across this additional session and the main experiment (yielding 900 sec total). Tuning during each attention condition was estimated separately.

Tuning changes for unattended categories were measured, and the direction of tuning shifts with respect to the attended categories was then inferred from these measurements. For this purpose, we used the semantic space that assesses of the similarity between the attended categories and remaining ones in terms of cortical representation. Specifically, if a single voxel’s tuning shifts towards categories similar to ‘humans’, then we should find that its tuning vector is closer to the ‘humans’ template than the ‘vehicles’ template in the semantic space. To test this prediction, voxel-wise tuning vectors in the control analysis and the template vectors for the attended categories were projected into the semantic space. Thereafter TSI was quantified following procedures in the main analysis, but to increase sensitivity we first computed the tuning change between the two attention conditions. We then computed an idealized tuning change between the template vectors in the semantic space. Finally, TSI was taken as the correlation between the actual and idealized tuning changes. As in the main analysis, tuning shifts toward the attended category will yield positive TSI values, whereas tuning shifts away from the target will yield negative TSI values.

The mean TSI is significantly greater than 0 in all subjects (Wilcoxon signed-rank test,  $p < 10^{-6}$ ; Supplementary Fig. 8). This result clearly shows that attention shifts tuning of unattended categories towards the attended category even when the targets are not present. Furthermore, the tuning shifts are in consistent directions across the main (Supplementary Fig. 7) and control analyses for an average of  $65.42 \pm 7.73\%$  (mean  $\pm$  s.d., averaged across subjects) of cortical voxels. While the direction of tuning shifts is highly consistent, the mean TSI is larger in the main analysis (when targets were present) than in the control analysis (when targets were excluded). TSI distributions are also less bimodal in the main analysis than in the control analysis. These differences are caused by two factors. First, the attentional effects on BOLD responses are strongest for the attended categories. Therefore, tuning changes obtained when the targets are present (main analysis) are naturally stronger than the tuning changes that occur when the targets are absent (control analysis). This reduces the TSI values in the control analysis.

Second, different metric spaces were used to estimate the TSI distributions in the main and the control analyses. In the main analysis, TSI was computed across 935 dimensions of the category model, and each category was treated as a separate dimension. As such, a tuning

shift in the direction of ‘humans’ represents tuning changes for ‘humans’ alone. Therefore, the main analysis only considers tuning changes for the attended categories, which account for  $38.79 \pm 0.07\%$  (mean $\pm$ s.d.) of tuning changes in cortical voxels. However, in the control analysis, TSI was computed across 6-dimensions of the semantic space that organizes categories according to semantic similarity. A tuning shift in the direction of ‘humans’ represents tuning changes for both ‘humans’ and nearby categories in this space. Therefore, the control analysis considers tuning changes for both attended and unattended categories, and tuning shifts toward the attended categories account for  $72.70 \pm 0.04\%$  (mean $\pm$ s.d.) of tuning changes in cortical voxels. This causes TSI distributions to be more bimodal in the control analysis.

### Cortical Flat Map Visualization

The cortical surface of each hemisphere was flattened after five relaxation cuts were applied to reduce distortions. For surface-based visualization, functional data were aligned to the anatomical data using in-house Matlab scripts (MathWorks Inc., Natick, MA). The functional data were then projected onto the cortical surface. Each point in the generated flat maps corresponded to an individual voxel.

A custom color map was designed to simultaneously visualize the cortical distribution of tuning strength for the attended categories. The tuning strengths (i.e.,  $s_H$  for ‘humans’ and  $s_V$  for ‘vehicles’) were measured as the correlations between the voxel-wise tuning vectors and the idealized templates tuned solely to these attended categories. Distinct colors were assigned to 6 landmark values of the pair ( $s_V$ ,  $s_H$ ): red for (0.75, 0), turquoise for (−0.75, 0), green for (0, 0.75), magenta for (0, −0.75), gray for (0, 0), and black for (−0.75, −0.75). The colors for the remaining values were linearly interpolated from these landmarks. A gray color was assigned to voxels with insignificant model weights.

A separate color map was designed to visualize the cortical distribution of semantic tuning. For this purpose, voxel-wise tuning vectors for each attention condition were projected into the semantic space. The first four PCs that captured approximately 20% of the variance were selected. The first PC mainly distinguishes categories with high versus low stimulus energy and so was not visualized. The projections onto the second, third and fourth PCs were assigned to the red, green and blue channels. Voxels with similar semantic tuning project to nearby points in the semantic space and so they were assigned similar colors. In this color map, voxels tuned for humans and communication verbs appeared in shades of green-cyan. Voxels tuned for animals and body parts appeared in yellow-green, whereas those tuned for movement verbs appeared in red. Voxels tuned for locations, roads, devices and artifacts appeared in shades of purple, whereas those tuned for buildings and furniture appeared in blue. Finally, voxels tuned for vehicles appeared in magenta. A gray color was assigned to voxels with insignificant model weights.

### Statistical Procedures

Statistical comparisons of prediction scores were based on raw correlation coefficients between the predicted and actual responses. Prediction scores were Fisher transformed; and one-sided t-tests were applied to assess significance. While this procedure is appropriate for

significance testing, noise in the measured BOLD responses biases raw correlation values downward<sup>30</sup>. Thus, to attain reliable estimates of model performance across subjects, correlation values were corrected for noise bias<sup>11</sup>.

Unless otherwise noted, all other comparisons were performed using one-sided non-parametric Wilcoxon signed-rank tests. All statistical significance levels were corrected for multiple comparisons using false-discovery-rate control<sup>23</sup>.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

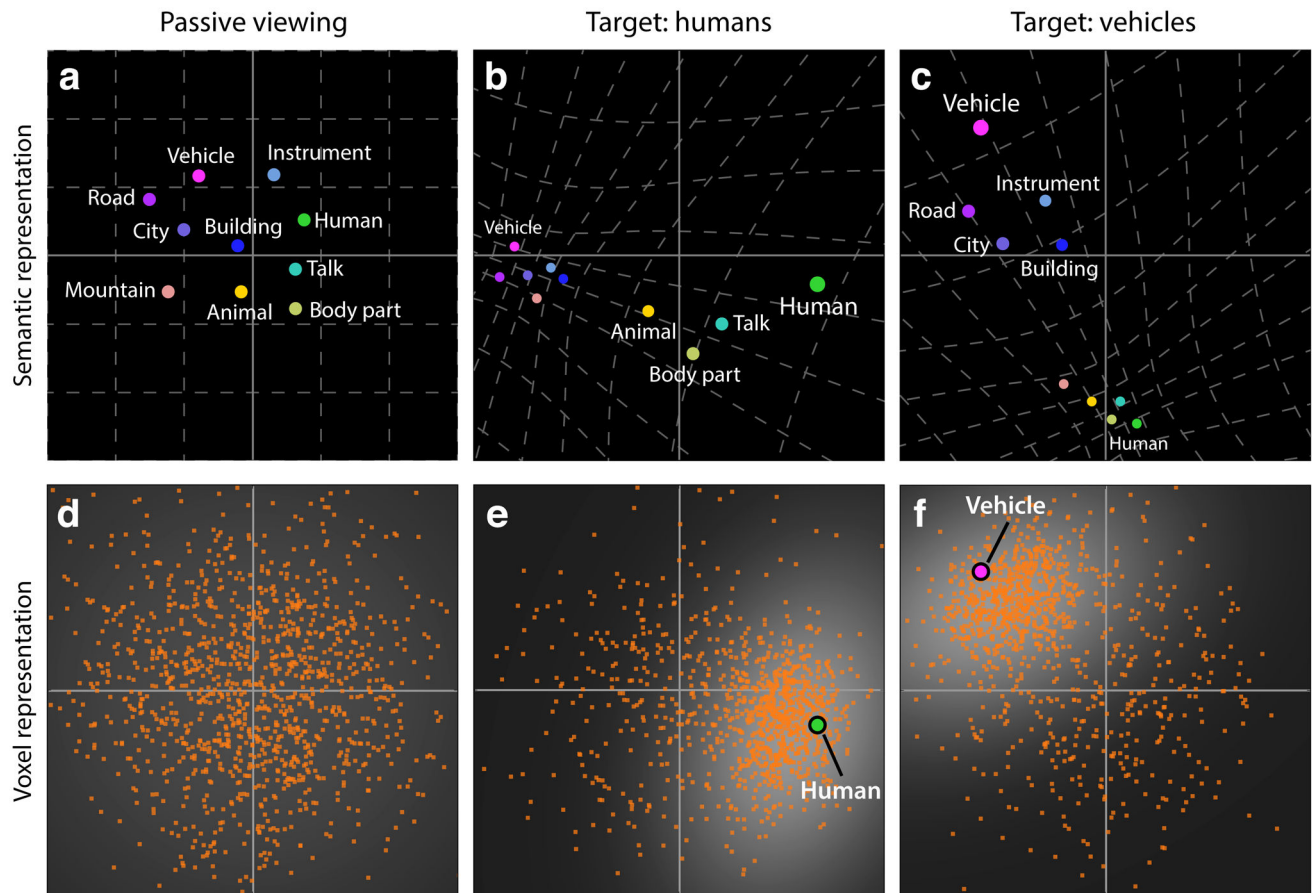
## Acknowledgments

This work was supported by the National Eye Institute (EY019684 and EY022454), and the Center for Science of Information (CSOI), an NSF Science and Technology Center, under grant agreement CCF-0939370. We thank David Whitney for discussions regarding this manuscript. We also thank J. Gao, N. Bilenko, T. Naselaris, A. Vu, and M. Oliver for their help in various aspects of this research.

## References

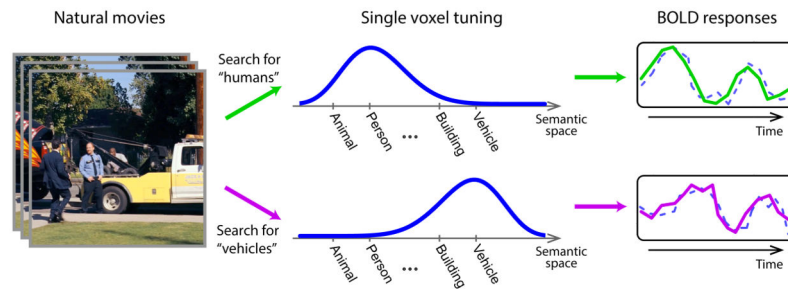
1. Olshausen BA, Anderson CH, Van Essen DC. A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *J Neurosci.* 1993; 13:4700–4719. [PubMed: 8229193]
2. Luck SJ, Chelazzi L, Hillyard SA, Desimone R. Neural mechanisms of spatial selective attention in areas V1, V2, and V4 of macaque visual cortex. *J Neurophysiol.* 1997; 77:24–42. [PubMed: 9120566]
3. McAdams CJ, Maunsell JH. Effects of attention on orientation-tuning functions of single neurons in macaque cortical area V4. *J Neurosci.* 1999; 19:431–441. [PubMed: 9870971]
4. Reynolds JH, Pasternak T, Desimone R. Attention increases sensitivity of V4 neurons. *Neuron.* 2000; 26:703–714. [PubMed: 10896165]
5. David SV, Hayden BY, Mazer JA, Gallant JL. Attention to stimulus features shifts spectral tuning of V4 neurons during natural vision. *Neuron.* 2008; 59:509–521. [PubMed: 18701075]
6. Connor CE, Preddie DC, Gallant JL, Van Essen DC. Spatial attention effects in macaque area V4. *J Neurosci.* 1997; 17:3201–3214. [PubMed: 9096154]
7. Asaad WF, Rainer G, Miller EK. Task-specific neural activity in the primate prefrontal cortex. *J Neurophysiol.* 2000; 84:451–459. [PubMed: 10899218]
8. Warden MR, Miller EK. Task-dependent changes in short-term memory in the prefrontal cortex. *J Neurosci.* 2010; 30:15801–15810. [PubMed: 21106819]
9. Johnston K, Everling S. Neural activity in monkey prefrontal cortex is modulated by task context and behavioral instruction during delayed-match-to-sample and conditional prosaccade-antisaccade tasks. *J Cognitive Neurosci.* 2006; 18:749–765.
10. Mazer JA, Gallant JL. Goal-related activity in V4 during free viewing visual search. Evidence for a ventral stream visual salience map. *Neuron.* 2003; 40:1241–1250. [PubMed: 14687556]
11. Huth AG, Nishimoto S, Vu AT, Gallant JL. A Continuous Semantic Space Describes the Representation of Thousands of Object and Action Categories across the Human Brain. *Neuron.* 2012; 76:1210–1224. [PubMed: 23259955]
12. Peelen MV, Fei-Fei L, Kastner S. Neural mechanisms of rapid natural scene categorization in human visual cortex. *Nature.* 2009; 460:94–97. [PubMed: 19506558]
13. Li FF, VanRullen R, Koch C, Perona P. Rapid natural scene categorization in the near absence of attention. *Proc Natl Acad Sci USA.* 2002; 99:9596–9601. [PubMed: 12077298]
14. O'Craven KM, Downing PE, Kanwisher N. fMRI evidence for objects as the units of attentional selection. *Nature.* 1999; 401:584–587. [PubMed: 10524624]

15. Reddy L, Kanwisher N. Category selectivity in the ventral visual pathway confers robustness to clutter and diverted attention. *Curr Biol*. 2007; 17:2067–2072. [PubMed: 17997310]
16. Bartels A, Zeki S. Functional brain mapping during free viewing of natural scenes. *Hum Brain Mapp*. 2004; 21:75–85. [PubMed: 14755595]
17. Hasson U, Nir Y, Levy I, Fuhrmann G, Malach R. Intersubject synchronization of cortical activity during natural vision. *Science*. 2004; 303:1634–1640. [PubMed: 15016991]
18. Kay KN, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008; 452:352–355. [PubMed: 18322462]
19. Nishimoto S, et al. Reconstructing visual experiences from brain activity evoked by natural movies. *Curr Biol*. 2011; 21:1641–1646. [PubMed: 21945275]
20. Naselaris T, Prenger RJ, Kay KN, Oliver M, Gallant JL. Bayesian reconstruction of natural images from human brain activity. *Neuron*. 2009; 63:902–915. [PubMed: 19778517]
21. Mitchell TM, et al. Predicting Human Brain Activity Associated with the Meanings of Nouns. *Science*. 2008; 320:1191–1195. [PubMed: 18511683]
22. Miller G. WordNet: a lexical database for English. *Commun ACM*. 1995; 38:39–41.
23. Benjamini Y, Yekutieli D. The control of the false discovery rate in multiple testing under dependency. *Ann Stat*. 2001; 29:1165–1188.
24. Corbetta M, Shulman GL. Control of goal-directed and stimulus-driven attention in the brain. *Nat Rev Neurosci*. 2002; 3:215–229.
25. Carter CS. Anterior Cingulate Cortex, Error Detection, and the Online Monitoring of Performance. *Science*. 1998; 280:747–749. [PubMed: 9563953]
26. Womelsdorf T, Anton-Erxleben K, Pieper F, Treue S. Dynamic shifts of visual receptive fields in cortical area MT by spatial attention. *Nat Neurosci*. 2006; 9:1156–1160. [PubMed: 16906153]
27. Haxby JV, et al. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*. 2001; 293:2425–2430. [PubMed: 11577229]
28. Turin G. An introduction to matched filters. *Information Theory, IRE Transactions on*. 1960; 6:311–329.
29. Navalpakkam V, Itti L. Search goal tunes visual features optimally. *Neuron*. 2007; 53:605–617. [PubMed: 17296560]
30. David SV, Gallant JL. Predicting neuronal responses during natural vision. *Network*. 2005; 16:239–260. [PubMed: 16411498]



**Figure 1. Tuning-shift hypothesis predicts that attention warps semantic representation**

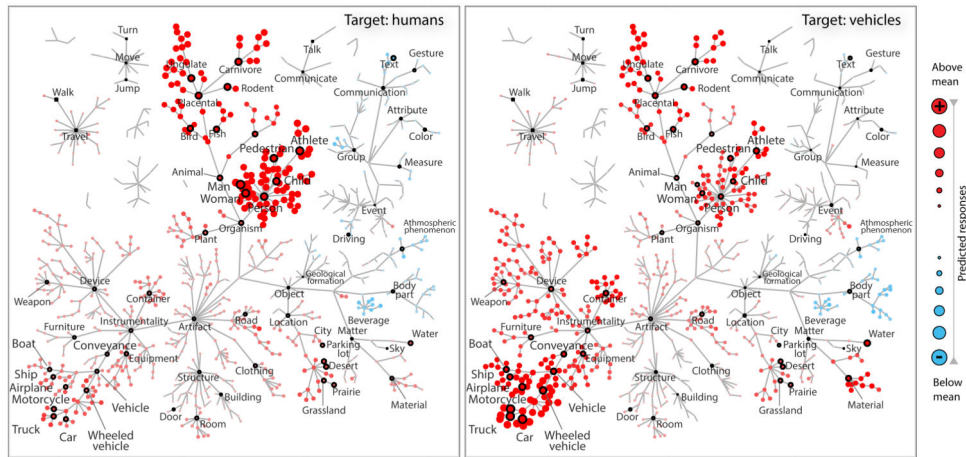
Hypothesized changes in semantic representation (top panel). Previous studies suggest that the brain represents categories by organizing them into a continuous space according to semantic similarity. **a**, During passive viewing semantically similar categories project to nearby points in the semantic space. **b–c**, The tuning-shift hypothesis predicts that attention to one specific category expands the representation of both the attended and nearby categories within the semantic space, and compresses the representation of distant categories. Attentional warping of semantic representation implies corresponding changes in voxel-wise semantic tuning (bottom panel). **d**, During passive viewing cortical voxels (orange dots) are tuned for different categories, and so they can also be visualized within the semantic space as in **a**. **e–f**, During visual search many voxels should shift their tuning toward the attended category in order to expand representation of the corresponding part of semantic space. This causes fewer voxels to be tuned for distant categories.



**Figure 2. Voxel-wise tuning vectors are measured from BOLD responses evoked by natural movies**

Tuning changes in single voxels are a unique, diagnostic aspect of the tuning-shift hypothesis. Therefore, to test this hypothesis we measured changes in voxel tuning during covert visual search for either ‘humans’ or ‘vehicles’ in complex natural movies. A separate category model was fit to each voxel within each attention condition, in order to optimally predict evoked BOLD responses (predicted response: dashed lines, measured response: solid lines). The category model gives voxel tuning under each condition, and tuning shifts can be identified by comparing tuning across conditions.

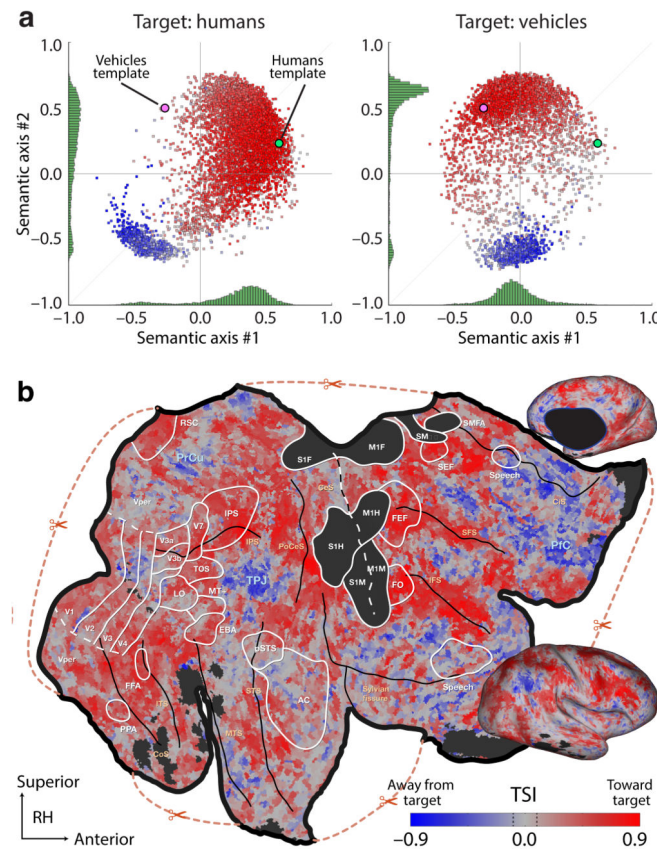




**Figure 3. Attentional tuning changes for a single voxel in LO**

Tuning for 935 object and action categories in a single voxel selected from lateral occipital complex (LO) in subject S1, during search for ‘humans’ (left) and for ‘vehicles’ (right). Each node in these graphs represents a distinct object or action, and a subset of the nodes has been labeled to orient the reader. The nodes have been organized using the hierarchical relations found in the WordNet lexicon. Red versus blue nodes correspond to categories that evoke above- and below-mean responses. The size of each node shows the magnitude of the category response (see legend on the right). This well-modeled LO voxel (a prediction score of 0.401) exhibits significant tuning changes across attention conditions (t-test,  $p < 10^{-6}$ ). The voxel is strongly tuned for the attended category in both conditions. Furthermore, significant albeit weaker tuning is observed for the unattended categories.

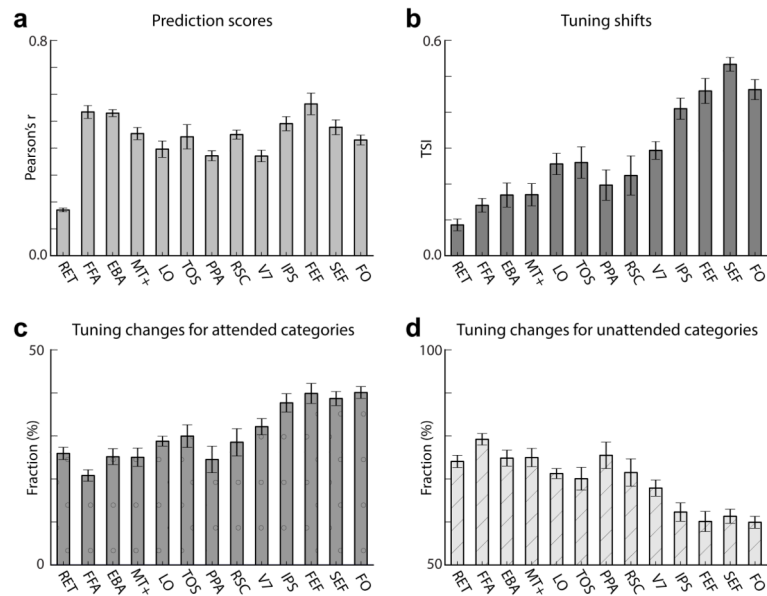




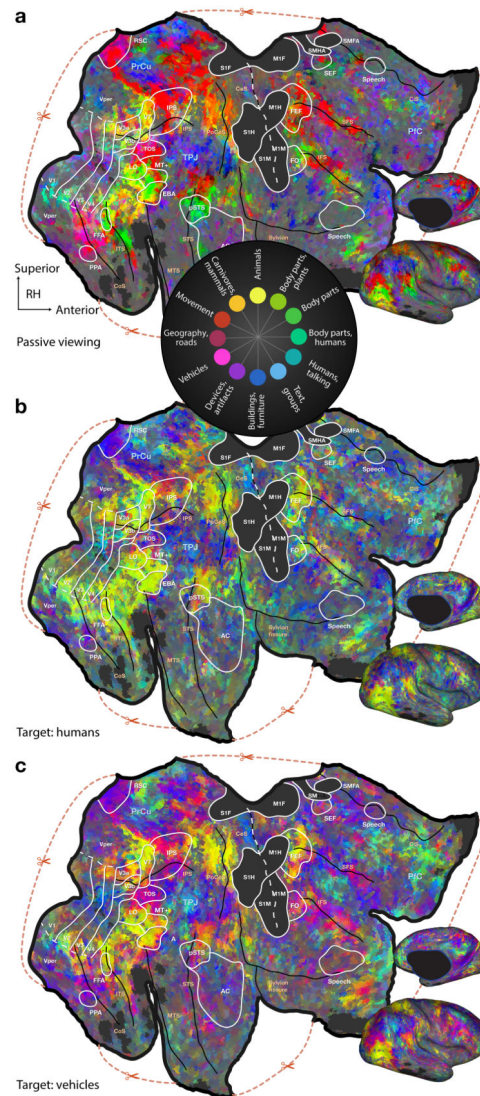
**Figure 4. Attention causes tuning shifts in single voxels**

**a**, Semantic tuning of single voxels during two attention conditions: search for ‘humans’ (left) and ‘vehicles’ (right). To assess attentional changes, voxel-wise tuning vectors were projected into a continuous semantic space. The semantic space was derived from principal components analysis (PCA) of tuning vectors measured during a separate passive-viewing task. Horizontal and vertical axes correspond to the second and third PCs (the first PC distinguishes categories with high versus low stimulus energy and so is not shown here). A total of 7785 well-modeled voxels with significant model weights (t-test,  $p < 0.05$ ) and high prediction scores (above mean plus 1 standard deviation) are shown for subject S1. Each voxel is represented with a dot whose color indicates the tuning shift index (TSI), red/blue for shifts toward/away from the target. The positions of the idealized templates for attended categories are shown in colored circles. The marginal distributions are displayed with separate histograms (green). Most well-modeled voxels strongly shift toward the attended category (Wilcoxon signed-rank test,  $p < 0.05$ ). **b**, The TSIs for subject S1 are shown on a cortical flat map of the right hemisphere. The color bar represents the 95% central range of TSIs and voxels with insignificant TSIs appear in gray ( $p > 0.05$ , within dashed black lines). Regions of fMRI signal dropout and motor areas excluded from all analyses are shown with dark gray patches. The boundaries of cortical areas identified by standard localizers are indicated with solid (functionally-inferred) and dashed (anatomically-inferred) white lines. Major anatomical landmarks (blue font) and sulci (orange font and black lines) are also labeled (see Supplementary Tables 1 and 2 for abbreviations). Voxels in many brain regions

shift their tuning toward the attended category. These include most of ventral-temporal cortex; lateral-occipital and intraparietal sulci (LO and IPS); inferior and superior frontal sulci (IFS and SFS); and dorsal bank of the cingulate sulcus (CiS). In contrast, the precuneus (PrCu), temporo-parietal junction (TPJ), anterior prefrontal cortex (Pfc), and areas along the anterior CiS shift their tuning away from the search target. Note: readers may explore the datasets at <http://gallantlab.org/brainviewer/cukuretal2013>.



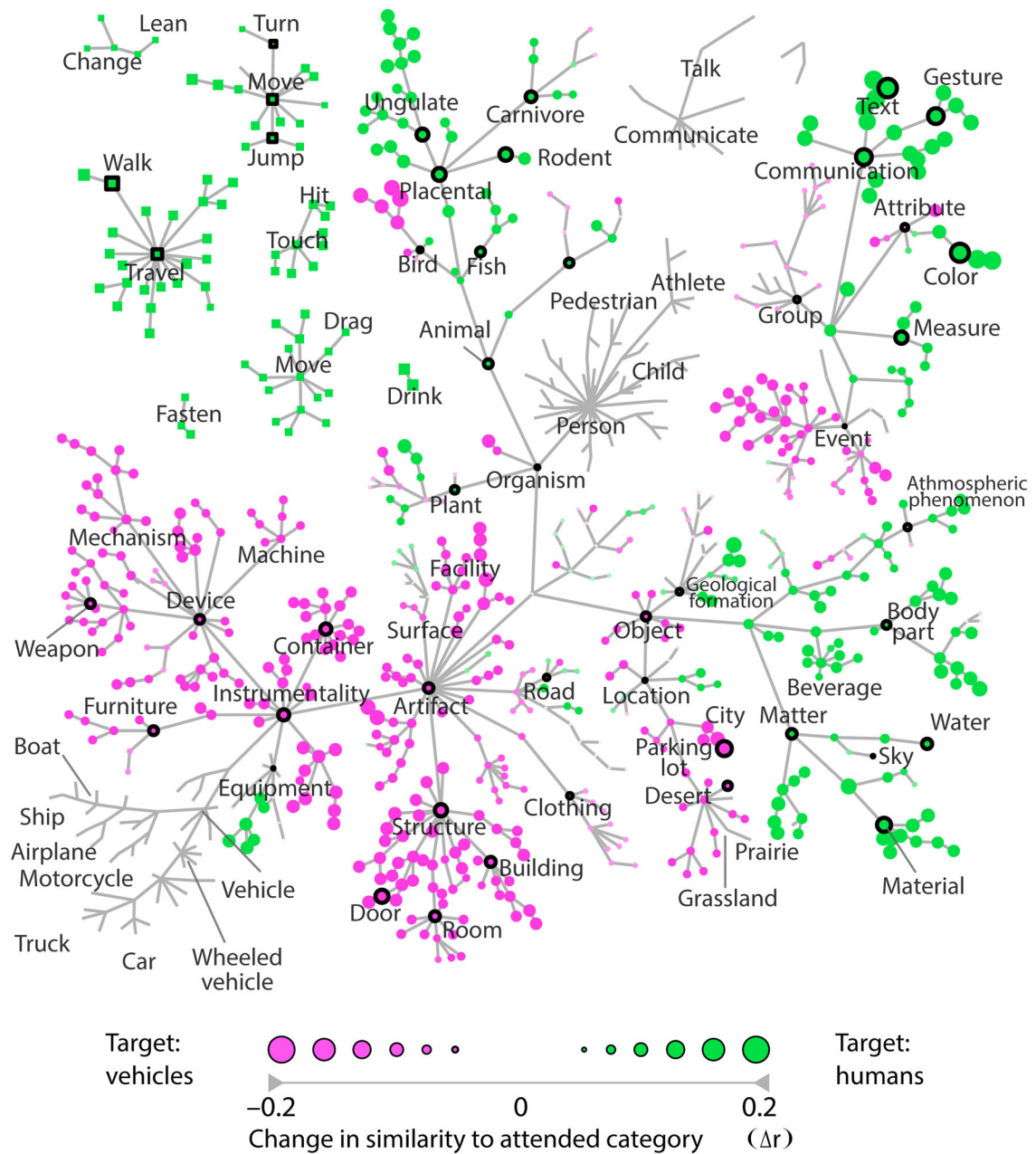
**Figure 5. Attention causes different degrees of tuning shifts in functional regions-of-interest**  
**a**, Prediction scores (Pearson's  $r$ ; mean $\pm$ s.e.m. results averaged across all 5 subjects). RET denotes the retinotopically-organized early visual areas V1–V3, and abbreviations for remaining regions-of-interest (ROIs) are listed in Supplementary Table 1. The average prediction score in category-selective areas in occipito-temporal cortex (FFA, EBA, LO, TOS) is  $0.48\pm0.07$  (mean $\pm$ s.d.); and the average prediction score in more anterior brain areas in frontal cortex (FEF, SEF, and FO) is  $0.49\pm0.07$  (mean $\pm$ s.d.). **a**, Tuning shift indices (TSI; mean $\pm$ s.e.m.) within functional ROIs. TSI is significantly greater than 0 in all ROIs (Wilcoxon signed-rank test,  $p<10^{-6}$ ). Furthermore, TSI increases towards relatively later stages of visual processing. **c**, Fraction of the overall tuning change (mean $\pm$ s.e.m.) explained by tuning changes for attended categories. **d**, Fraction of the overall tuning change (mean $\pm$ s.e.m.) explained by tuning changes for unattended categories (i.e., excluding both 'humans' and 'vehicles'). The degree of tuning shift (i.e., TSI) is positively correlated with the fraction of variance explained by tuning changes for attended categories ( $r=0.86\pm0.02$ , t-test,  $p<10^{-6}$ ).



**Figure 6. Semantic tuning for unattended categories shifts toward the attended category even when no targets are present**

**a**, Distribution of semantic tuning across the cortex (subject S1, right hemisphere) during passive viewing. Tuning was estimated from responses to all available movie clips. A four-dimensional semantic space was derived from these data using PCA. The tuning vector for each cortical voxel was then projected into this space; and the projections onto the second, third and fourth PCs were assigned to the red, green and blue channels. Here voxels with similar tuning project to nearby points in the semantic space and so they are assigned similar colors (see legend). Insignificant voxels are shown in gray. Yellow-green voxels are more selectively tuned for animals and body parts, and purple-red voxels are more selectively tuned for geographic locations and movement. Anatomical landmarks are labeled as in Fig. 4b. **b**, Distribution of semantic tuning for the same subject as in panel (a), but during search for ‘humans’. Tuning was estimated only from responses evoked by movie clips in which the target did not appear. Color assignment same as in panel (a). Yellow-green voxels that are tuned for animals and body parts predominate during search for ‘humans’. Many voxels

in posterior areas that are tuned for vehicles under passive viewing (e.g., PPA, RSC, and TOS) shift their tuning away from vehicles; and many voxels that are not tuned for humans under passive viewing (in FEF, FO, IPS, PfC and insular cortex) shift their tuning toward humans. **c**, Distribution of semantic tuning for the same subject as in panel (a), but during search for ‘vehicles’. Tuning was estimated only from responses evoked by movie clips in which the target did not appear. Color assignment same as in panel (a). Purple-magenta voxels that are tuned for geographic locations and artifacts predominate during search for ‘vehicles’. Many voxels in posterior areas that are tuned for humans under passive viewing (e.g., EBA, FFA, TPJ, and PrCu) shift their tuning away from humans; and many voxels that are not tuned for vehicles under passive viewing (in FEF, FO, IPS, PfC and insular cortex) shift their tuning toward vehicles. Note: readers may explore the datasets at <http://gallantlab.org/brainviewer/cukuretal2013>.



The tuning-shift hypothesis predicts that attention expands the representation of unattended categories that are nearby the attended category within the semantic space. This implies that the representation of unattended categories that are semantically similar to the target will shift toward the representation of the attended category. To address this issue we measured the similarity of BOLD-response patterns evoked by unattended categories to those evoked by the attended category. In each subject, response patterns were estimated across a total of 4245–7785 well-modeled voxels that were used in the main analysis. The response patterns for unattended and attended categories were estimated using target-absent and target-present



movie segments, respectively. The similarity of response patterns was quantified using Pearson's correlation ( $r$ ), and the results were averaged across subjects. In this figure each node represents a distinct object or action, and some nodes have been labeled to orient the reader. The nodes have been organized using the hierarchical relations found in the WordNet lexicon. The size of each node shows the magnitude of change in similarity (Wilcoxon signed-rank test,  $p < 10^{-4}$ ; see legend at the bottom). During search for 'humans', representations of semantically similar categories (e.g., animals, body parts, action verbs and natural materials) shift toward the representation of 'humans' (green nodes). During search for 'vehicles', representations of semantically similar categories (e.g., tools, devices, and structures) shift toward the representation of 'vehicles' (magenta nodes).